

Data Analysis Projects

MoBi 4. FS - SS2019



Medizinische Fakultät Heidelberg

Concept



Medizinische Fakultät Heidelberg

- **Project-oriented teaching:** provide hands-on experience with data analysis and programming
- **Goals**
 - experience data analysis challenges on real datasets related to research question
 - experience team work, also outside of your team!
 - learn to use modern data analysis tools: R / Python / markdown / notebooks / github

Research topics / projects



Medizinische Fakultät Heidelberg

- **5 research projects** have been defined
- Each topic has up to **5 sub-projects**
- Each project will be worked out by **groups of 4 students**
- **one supervisor and master tutor per project**
- **Role of the master tutors:**
 - weekly meetings with groups working on project
(Wednesday 10am-1pm)
- Meeting rooms: BioQuant SR42/43 + IPMB meeting room 5th floor
2 other meeting places must be found (lounge corner at BioQuant?...)
- Tutors:
Valentina Giunchiglia; Julia Rühle; David Schwarzenbacher; Nicolas Peschke; Alexander Mattausch

Timeline



Medizinische Fakultät Heidelberg

17/04

We do...

Presentation of the projects

24/04

Presentation of R markdown
and github

15/05

You do...

selection of projects
and teams;
registration

24/07

(25/07 for Project 03)



Presentation of
project proposal
(10 + 10 min)

Final presentation
(15+10 minutes)

Project proposal (15/05)



Medizinische Fakultät Heidelberg

- During the project proposal presentation, you should
 - review some of the references given in the project description
 - explain what the questions / challenges are
 - describe which of these questions you want to address in your project
 - indicate a approximate timeline
 - ▶ milestones = important steps in the analysis
 - ▶ when these milestones should be achieved
- Presentation in front of the project supervisors
 - 10 minutes presentation
 - 10 minutes discussion / questions
- *All team members are expected to contribute!*

Projects / sub-projects



Medizinische Fakultät Heidelberg

- **Project 01: *Genetic interactions in cancer***
(Ashwini Sharma / Carl Herrmann)
 - Data types: gene expression / gene mutations / gene knockdown / CNV
- **Project 02: *Cellular response to drug perturbation***
(Nicolas Palacio / Javier Perales)
 - Data types: gene expression treated / untreated / mutations / metadata
- **Project 03: *Biomedical image analysis***
(Karl Rohr / Christian Ritter)
 - Data types: MNIST images / cell nuclei images
- **Project 04: *Programming k-means***
(Thorsten Beier)
 - Data types: scRNA-seq
- **Project 05: *Cancer DNA Methylation***
(Matthias Schlesner / Christian Heyer)
 - Data types: DNA methylation WGBS

Project selection / registration

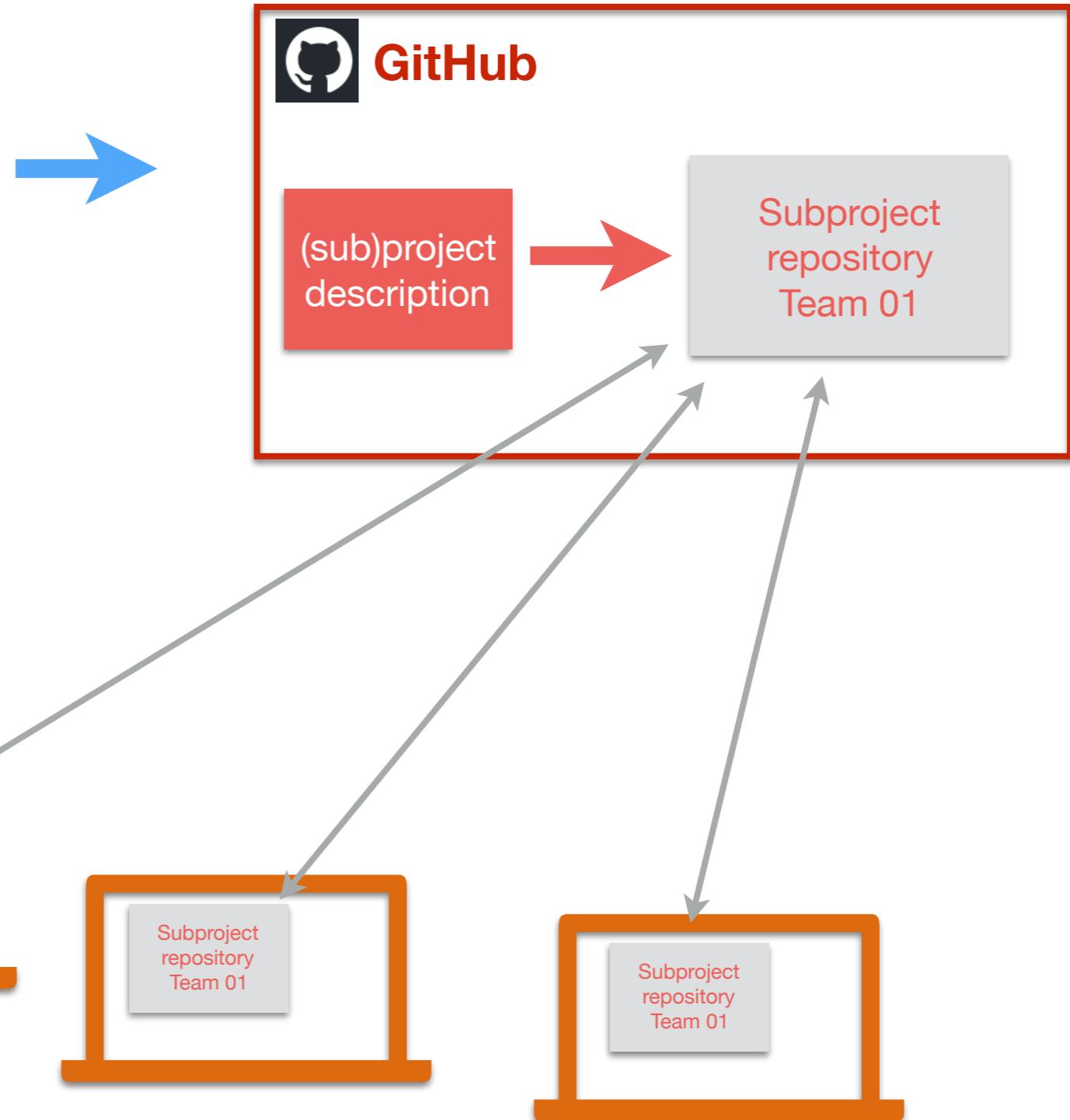


Medizinische Fakultät Heidelberg

- Listen to the description of the projects / sub-projects
- Check this webpage
<https://datascience-mobi.github.io/>
- Once you have selected your team and project, **register your team** in the Google sheet
[https://docs.google.com/spreadsheets/d/1LEQLH2LaDulMq3Qepx-7-5KWjgZjK4dfWDDIbm1vu0Q/edit?
usp=sharing](https://docs.google.com/spreadsheets/d/1LEQLH2LaDulMq3Qepx-7-5KWjgZjK4dfWDDIbm1vu0Q/edit?usp=sharing)
- ***Selection of the projects should be done by 24/04 10am !***

Organizing your work

Website
datascience-mobi.github.io

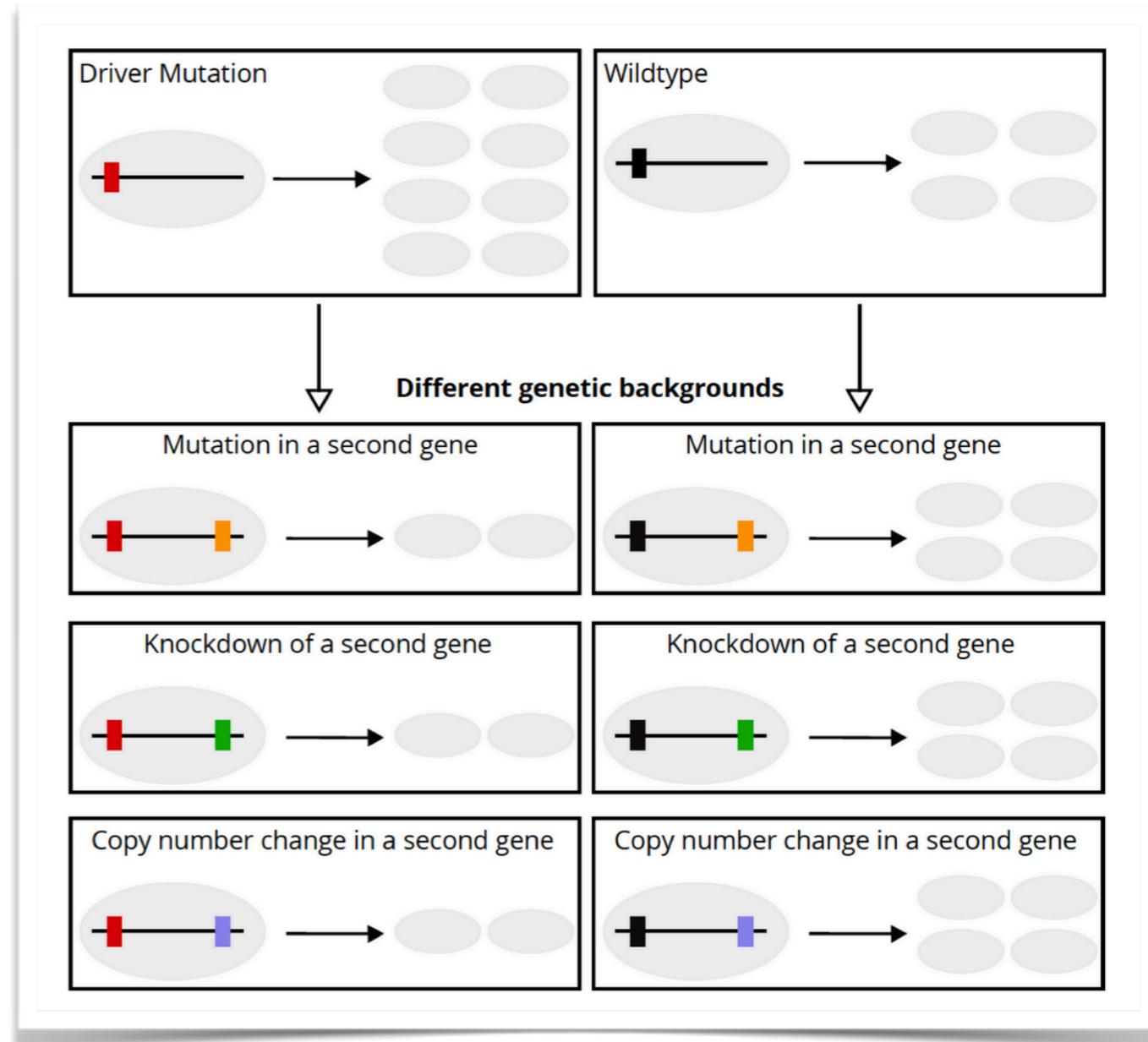


jupyter

R Markdown
from R Studio®

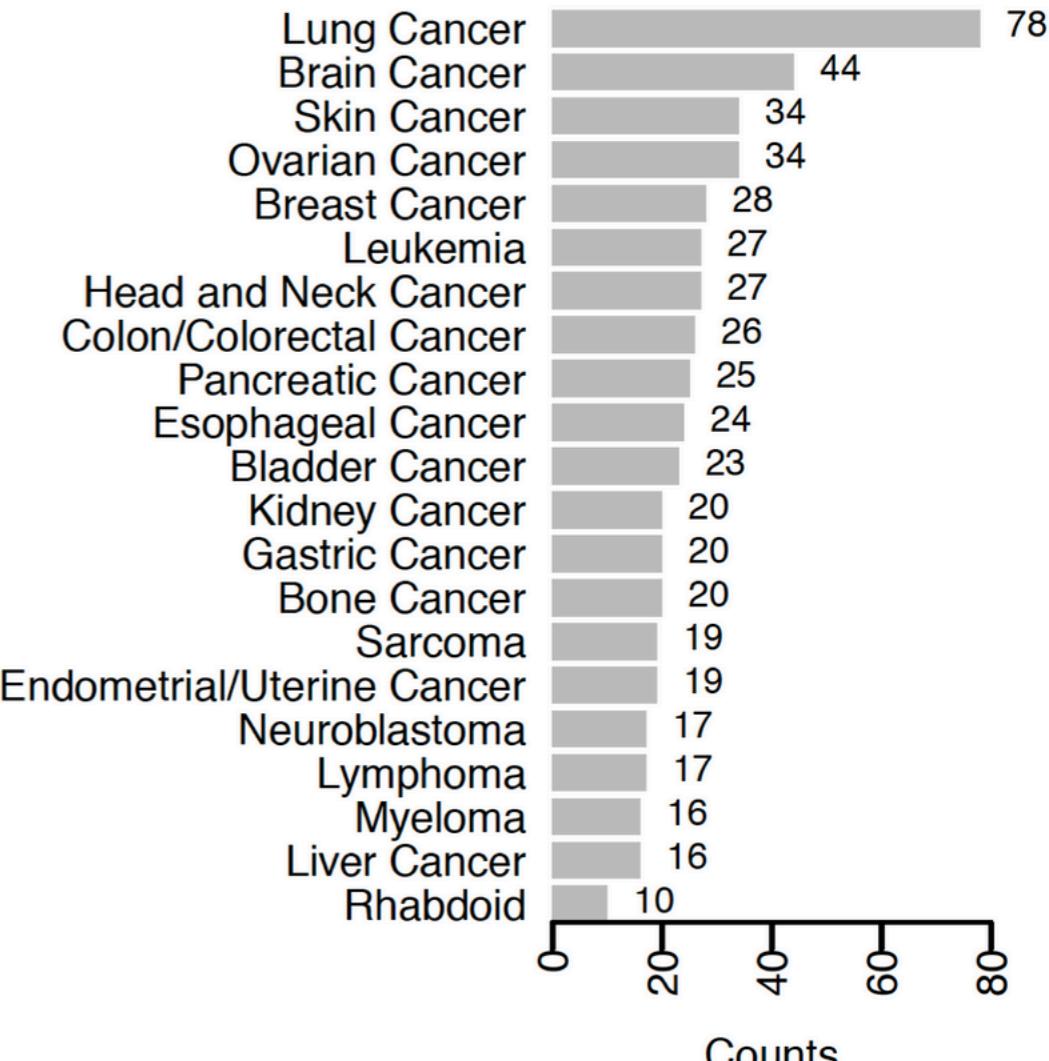
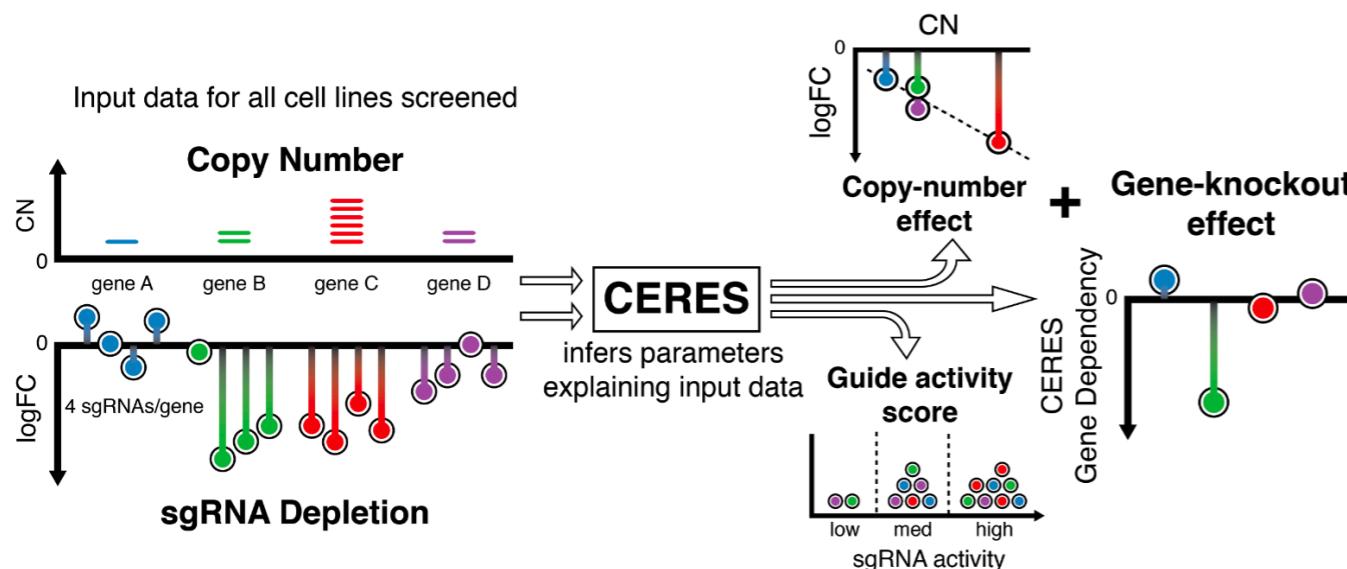
Project 01: genetic interactions in Cancer cells

- Find **synthetic lethal interactions** in cancer cell lines for driver mutations
- Which mutation / copy-number alteration does lead together with a driver mutation to a significant phenotype in cancer cells?
→ identify potential therapeutic targets



Project 01: genetic interactions in Cancer cells

- Knockdown assays (Project DepMap/CERES): CRISPR-Cas9 for ~17.000 genes across 544 cell lines



[Meyer et al., 2017]



Datasets

- **Gene expression** in various cancer cell lines
- **Mutations** in genes in each cell line & copy-number alterations in each cell line
- **Knock-down sensitivity** in all cell lines for 17.000 genes
- **Annotations** for cell lines (tissue, ...)

	ACH-00004	ACH-00005	ACH-00007	ACH-00009	ACH-00011	ACH-00012	ACH-00013
TSPAN6	2.61588707	3.06695024	4.06608919	6.50795317	4.57773093	5.82094479	5.53387478
TNMD	0.00000000	0.00000000	0.00000000	0.09761080	0.00000000	0.00000000	0.05658353
DPM1	5.32337007	5.76261458	5.88996020	7.98162436	5.53574194	6.00652245	7.53216108
SCYL3	2.40599236	2.99276843	3.04963077	2.24792751	2.08746284	1.91456452	1.91073266
C1orf112	3.90207358	5.35966172	3.76022095	4.49121176	2.64616266	3.75167795	3.54720302
FGR	0.92599942	0.23878686	0.02856915	0.00000000	0.00000000	0.02856915	0.02856915
CFH	4.88898672	5.70016225	0.01435529	0.02856915	0.35614381	2.68929916	1.98550043

Hugo_Symbol	Entrez_Gene_Id	NCBI_Build	Chromosome	Start_position	End_position	Strand	Variant_Classification
RNF207	388591	37	1	6279339	6279339	+	Missense_Mutation
PLEKHG5	57449	37	1	6533165	6533165	+	Missense_Mutation
PDPN	10630	37	1	13940848	13940848	+	Missense_Mutation
CASP9	842	37	1	15819484	15819484	+	Missense_Mutation
RAP1GAP	5909	37	1	21924552	21924552	+	Missense_Mutation
C1QC	714	37	1	22974054	22974054	+	Silent
CNKS1	10256	37	1	26514778	26514778	+	Frame_Shift_Del
AHDC1	27245	37	1	27875454	27875454	+	Missense_Mutation
COL16A1	1307	37	1	32133218	32133218	+	Silent
CSMD2	114784	37	1	32133218	34383844	+	Silent
LRRC8C	84230	37	1	90179033	90179033	+	Missense_Mutation
SPAG17	200162	37	1	118567972	118567973	+	Frame_Shift_Del
TBX15	6913	37	1	119427935	119427935	+	Missense_Mutation

	ACH-00004	ACH-00005	ACH-00007	ACH-00009
A1BG	0.1346453616	-0.212445060	0.043317923	0.0705119999
A1CF	0.0755362715	0.233123579	0.066837574	0.0084297636
A2M	-0.1402086015	0.044364933	-0.036196515	0.0271141959
A2ML1	0.0139284337	0.173837240	0.134781001	0.0559267727
A3GALT2	0.0291310328	-0.124389318	0.082995584	0.0463253889
A4GALT	-0.1472838445	-0.298849014	0.119084008	0.0159682666
A4GNT	0.2758291936	0.120259815	0.057116006	0.0535023006
AAAS	-0.3636329615	-0.339925280	-0.352541473	-0.4988600588

Project 01: genetic interactions in Cancer cells



Medizinische Fakultät Heidelberg

- Select cancer type and corresponding cell lines
- Determine driver mutations from the literature (e.g. EGFR mutations in lung cancer)
- Determine potential synthetic lethal mutations/ copy-number alterations from knock-down screens by splitting **mutated/non-mutated** cell lines

non-
mutated
cell lines

mutated
cell lines



Medizinische Fakultät Heidelberg

Brief intro to Git(Hub)

Git(Hub)



Medizinische Fakultät Heidelberg

- Git is a **version control system**:
 - allows simultaneous work of different people on the same project
 - tracks the changes ('**commits**') made by each member
 - helps solve the **conflicts** between various versions
- GitHub is a platform which hosts Git projects ('**repositories**')
 - is free to use
 - required to create a (free) account
 - can be used in command line or using GUI tools ('**GitHub Desktop**')

Git repository

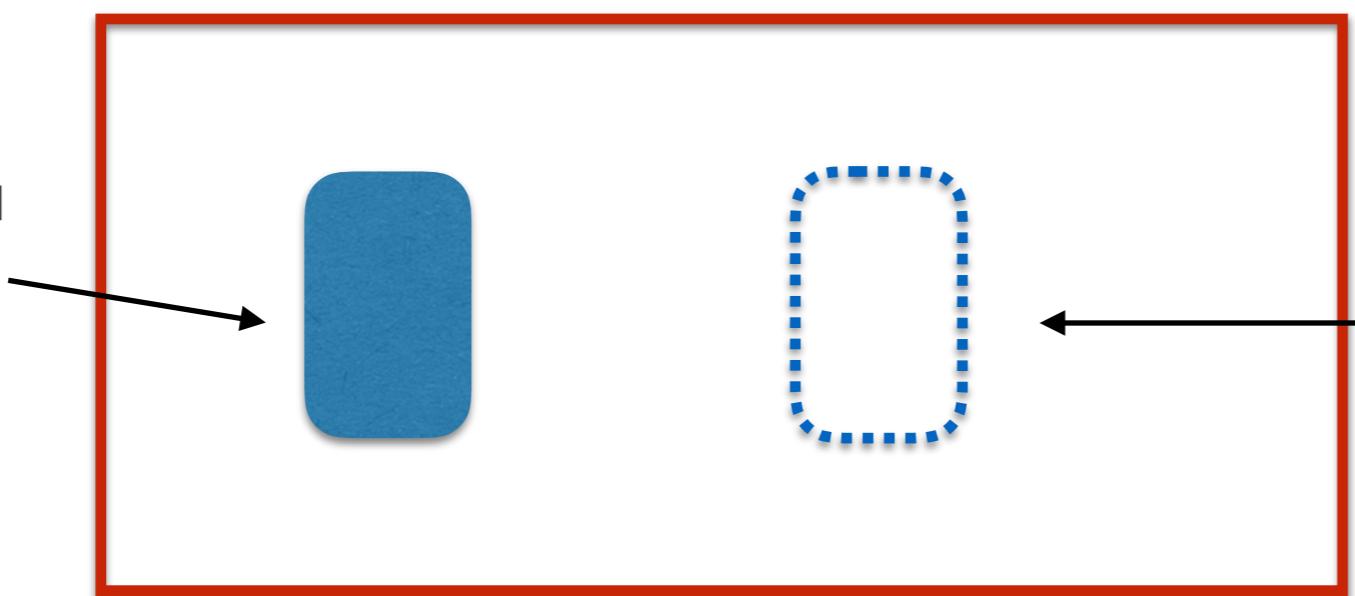


Medizinische Fakultät Heidelberg

repository



this file is registered
in the database
(‘committed’)



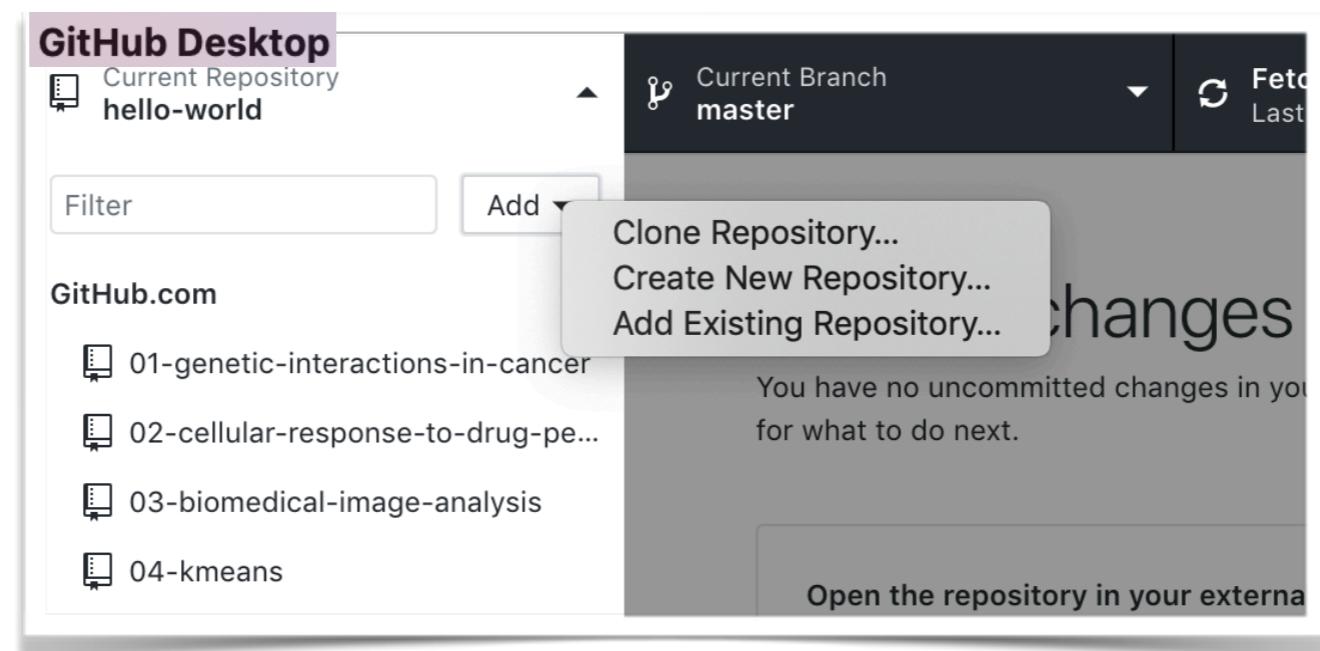
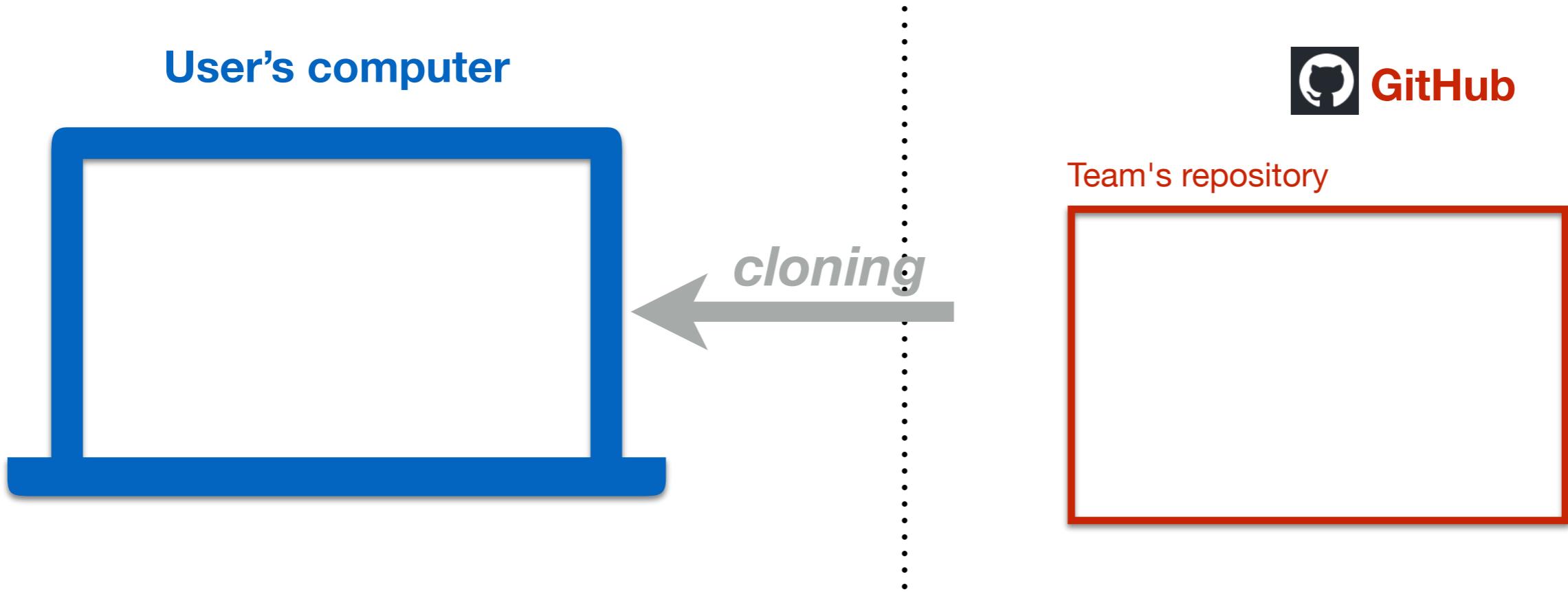
this file exists, but
has not been
registered yet
OR
the file has been modified, but
the changes have not been
registered yet

1. Cloning an existing repository



Medizinische Fakultät Heidelberg

User's computer





erg

Current Repository **hello-world**

Current Branch **master**

Fetch origin
Last fetched 19 minutes ago

Changes History

0 changed files

No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

Open the repository in your external editor
Configure which editor you wish to use in [preferences](#)

Repository menu or ⌘ ⌘ A

[Open in Visual Studio Code](#)

View the files in your repository in Finder
Repository menu or ⌘ ⌘ F

[Show in Finder](#)

Open the repository page on GitHub in your browser
Repository menu or ⌘ ⌘ G

[View on GitHub](#)

hd Summary (required)

Description

+

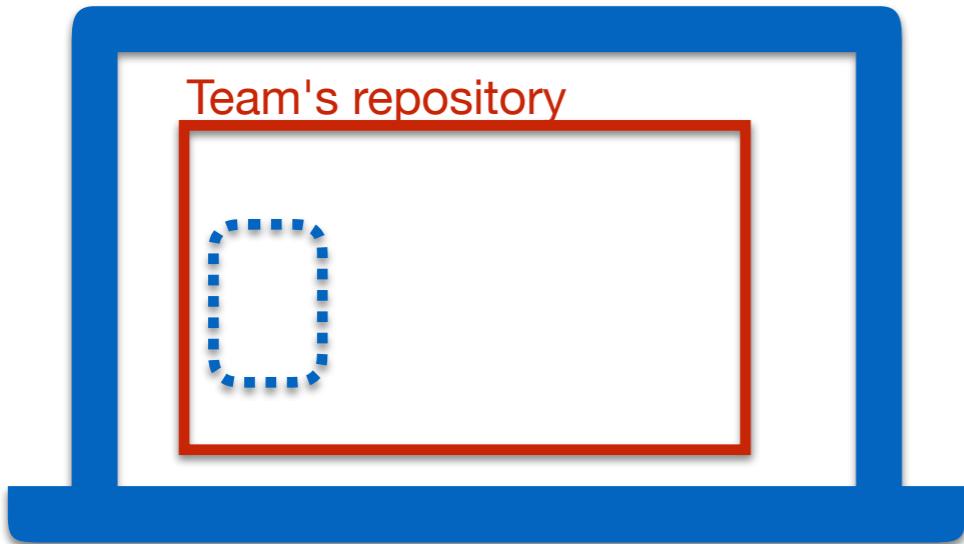
Commit to master

2. creating a local file



Medizinische Fakultät Heidelberg

User's computer



Team's repository

 GitHub



- When a new file is added / modified in the local folder, it is not yet registered in the git database!
- it first needs to be **committed**



lberg

Current Repository **hello-world** Current Branch **master** Fetch origin Last fetched 22 minutes ago

Changes 1 History my_markdown.Rmd +

1 changed file
 my_markdown.Rmd +

new file created locally

```
@@ -0,0 +1,30 @@
1 +---
2 +title: "My first markdown"
3 +author: "Carl Herrmann"
4 +date: "4/23/2019"
5 +output: html_document
6 +---
7 +
8 +```{r setup, include=FALSE}
9 +knitr::opts_chunk$set(echo = TRUE)
10 +```
11 +
12 +## R Markdown
13 +
14 +This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.
15 +
16 +When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:
17 +
18 +```{r cars}
19 +summary(cars)
20 +```
21 +
22 +## Including Plots
23 +
24 +You can also embed plots, for example:
```

hd su Create my_markdown.Rmd

Description

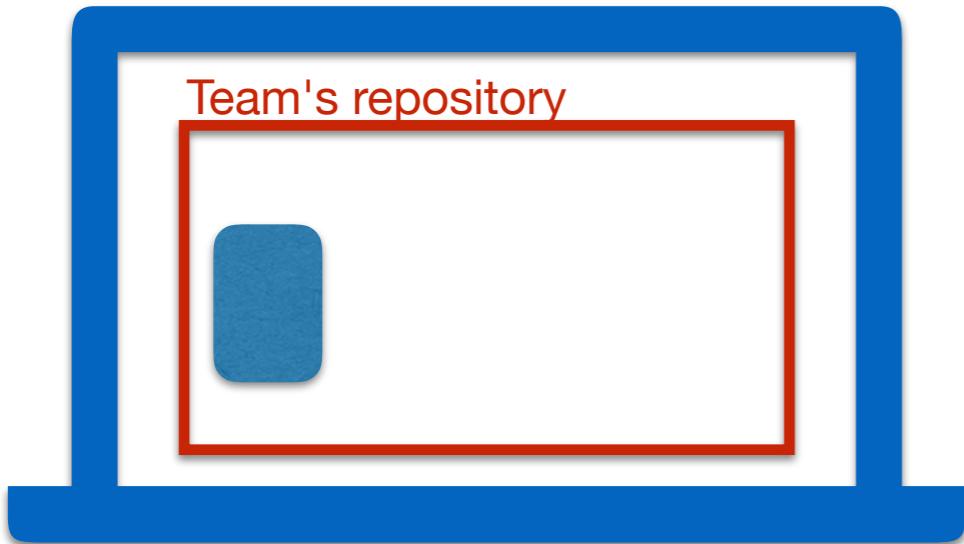
Commit to master

2. Adding a file



Medizinische Fakultät Heidelberg

User's computer



 GitHub

Team's repository



- When a new file is added / modified in the local folder, it is not yet registered in the git database!
- it first needs to be **committed**



lberg

Current Repository **hello-world** Current Branch **master** Fetch origin Last fetched 22 minutes ago

Changes 1 History my_markdown.Rmd +

1 changed file my_markdown.Rmd +

indicate the type of changes made and commit

Create my_markdown.Rmd

Description

Commit to master

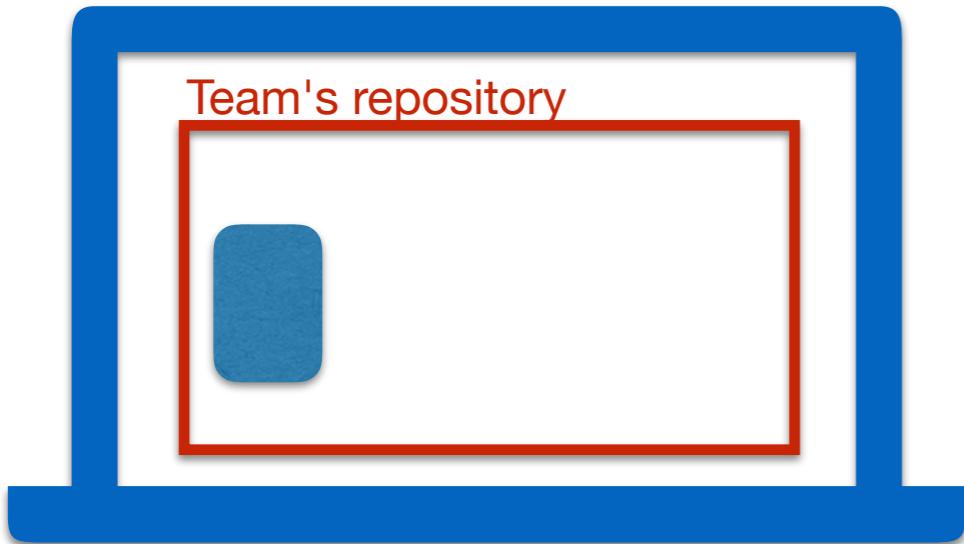
```
@@ -0,0 +1,30 @@
1 +---
2 +title: "My first markdown"
3 +author: "Carl Herrmann"
4 +date: "4/23/2019"
5 +output: html_document
6 +---
7 +
8 +```{r setup, include=FALSE}
9 +knitr::opts_chunk$set(echo = TRUE)
10 +```
11 +
12 +## R Markdown
13 +
14 +This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.
15 +
16 +When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:
17 +
18 +```{r cars}
19 +summary(cars)
20 +```
21 +
22 +## Including Plots
23 +
24 +You can also embed plots, for example:
```

2. Adding a file



Medizinische Fakultät Heidelberg

User's computer



Team's repository



- the file is now committed to the local git repository
- it needs to be pushed to the remote repository on GitHub



Current Repository **hello-world** Current Branch **master** Push origin Last fetched 30 minutes ago 1 ↑

Changes History 0 changed files

No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

Push 1 commit to the origin remote
You have one local commit waiting to be pushed to GitHub
Always available in the toolbar when there are local commits waiting to be pushed or ⌘ P

Push origin

Open the repository in your external editor
Configure which editor you wish to use in [preferences](#)
Repository menu or ⌘ ↑ A

Open in Visual Studio Code

View the files in your repository in Finder
Repository menu or ⌘ ↑ F

Show in Finder

Open the repository page on GitHub in your browser
Repository menu or ⌘ ↑ G

View on GitHub

hd Summary (required)

Description

+

Commit to **master**

Committed just now Create my_markdown.Rmd

Undo



Medizinische Fakultät Heidelberg

test repository

[Edit](#)

[Manage topics](#)

⌚ 2 commits ⚡ 1 branch ⚡ 0 releases ⚡ 1 contributor

Branch: master ▾ [New pull request](#) [Create new file](#) [Upload files](#) [Find File](#) [Clone or download ▾](#)

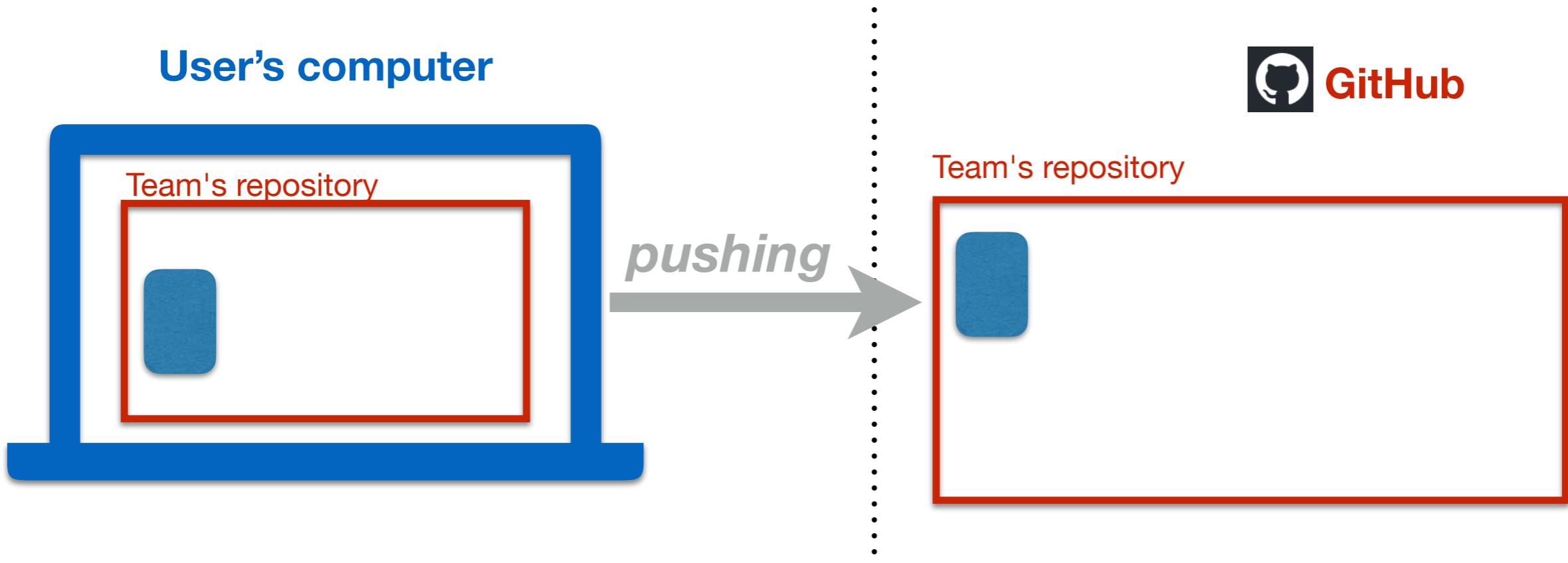
 carlherrmann	Create my_markdown.Rmd	Latest commit b2b0bdf 2 minutes ago
 README.md	Initial commit	37 minutes ago
 my_markdown.Rmd	Create my_markdown.Rmd	2 minutes ago

 README.md [!\[\]\(e551b45089a3dd3ec457d321817b7ac6_img.jpg\)](#)

hello-world

test repository

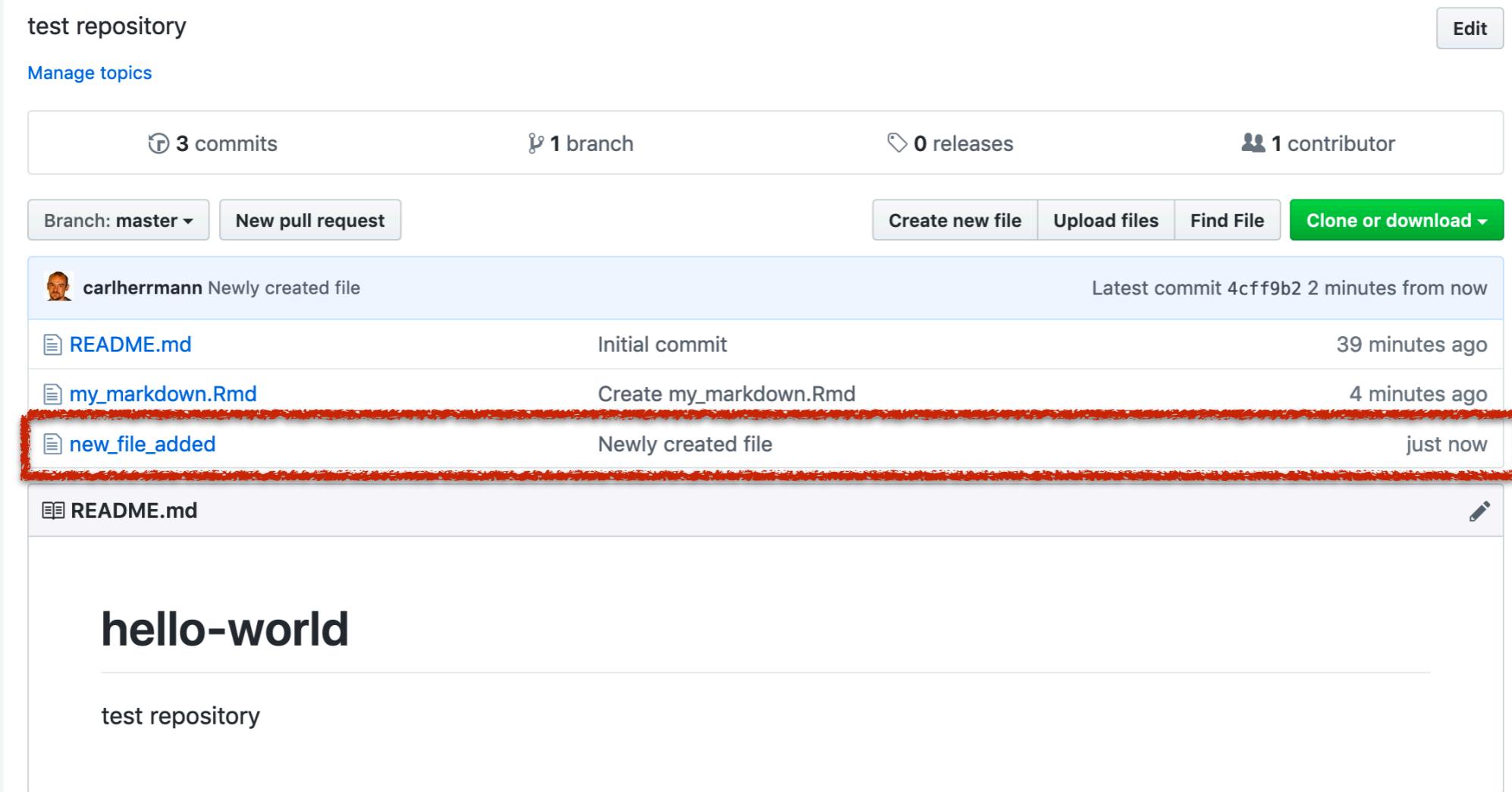
2. Adding a file



- the file is no committed to the local git repository
- it needs to be pushed to the remote repository on GitHub

3. Pulling from the remote repository

- Someone (probably one of your team mates) has added a new file into the remote repository
- It is not yet in your local repository and need to be **pulled**



The screenshot shows a GitHub repository named "test repository". The repository summary indicates 3 commits, 1 branch, 0 releases, and 1 contributor. The "Clone or download" button is highlighted in green. The commit history lists three entries:

- carlherrmann Newly created file (Latest commit 4cff9b2 2 minutes from now)
- Initial commit (39 minutes ago)
- Create my_markdown.Rmd (4 minutes ago)
- new_file_added** Newly created file (just now)

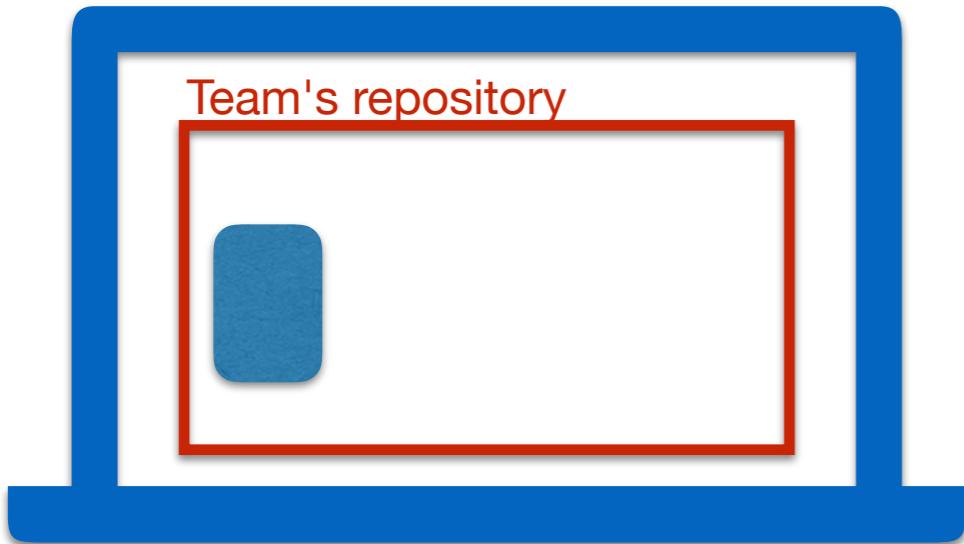
A red box highlights the most recent commit, "new_file_added". Below the commit history, there is a preview of the file "README.md" which contains the text "hello-world".

3. Pulling from the remote repository



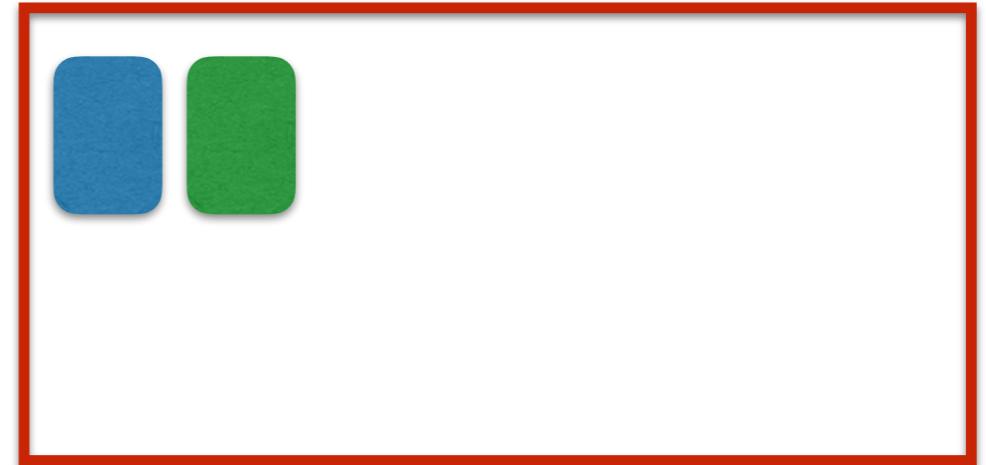
Medizinische Fakultät Heidelberg

User's computer



 GitHub

Team's repository





Current Repository **hello-world**

Current Branch **master**

Fetch origin
Last fetched 4 minutes ago

Changes History

0 changed files

No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

Open the repository in your external editor
Configure which editor you wish to use in [preferences](#)

Repository menu or ⌘ ⌘ A

Open in Visual Studio Code

View the files in your repository in Finder
Repository menu or ⌘ ⌘ F

Show in Finder

Open the repository page on GitHub in your browser
Repository menu or ⌘ ⌘ G

View on GitHub

hd su Summary (required)

Description

Commit to master



Current Repository **hello-world**

Current Branch **master**

Pull origin
Last fetched just now

Changes History

0 changed files

No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

Pull 1 commit from the origin remote
The current branch (master) has a commit on GitHub that does not exist on your machine.

Always available in the toolbar when there are remote changes or

Pull origin

Open the repository in your external editor
Configure which editor you wish to use in [preferences](#)

Repository menu or

Show in Visual Studio Code

View the files in your repository in Finder
Repository menu or

Show in Finder

Open the repository page on GitHub in your browser
Repository menu or

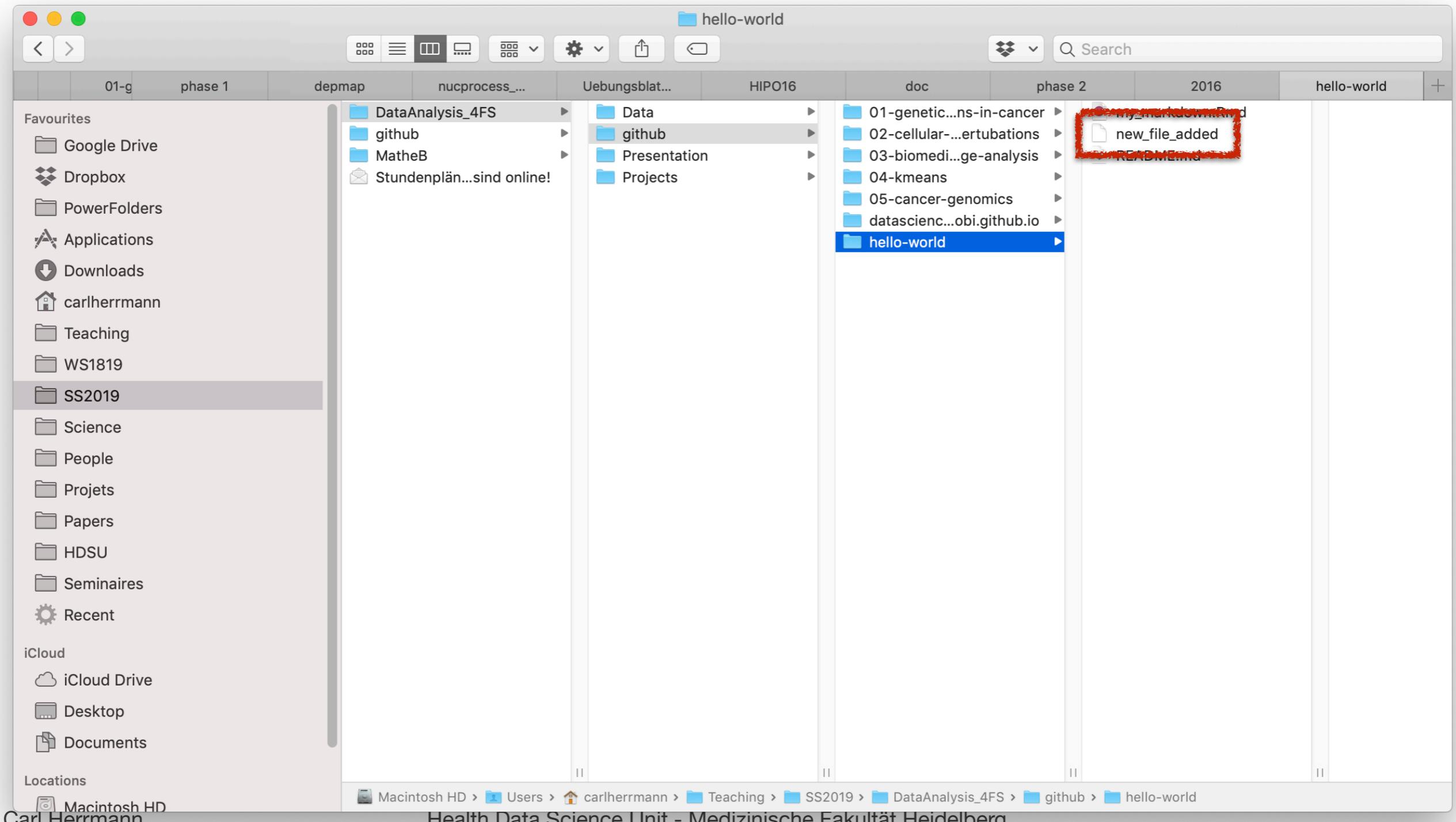
View on GitHub

hd Summary (required)

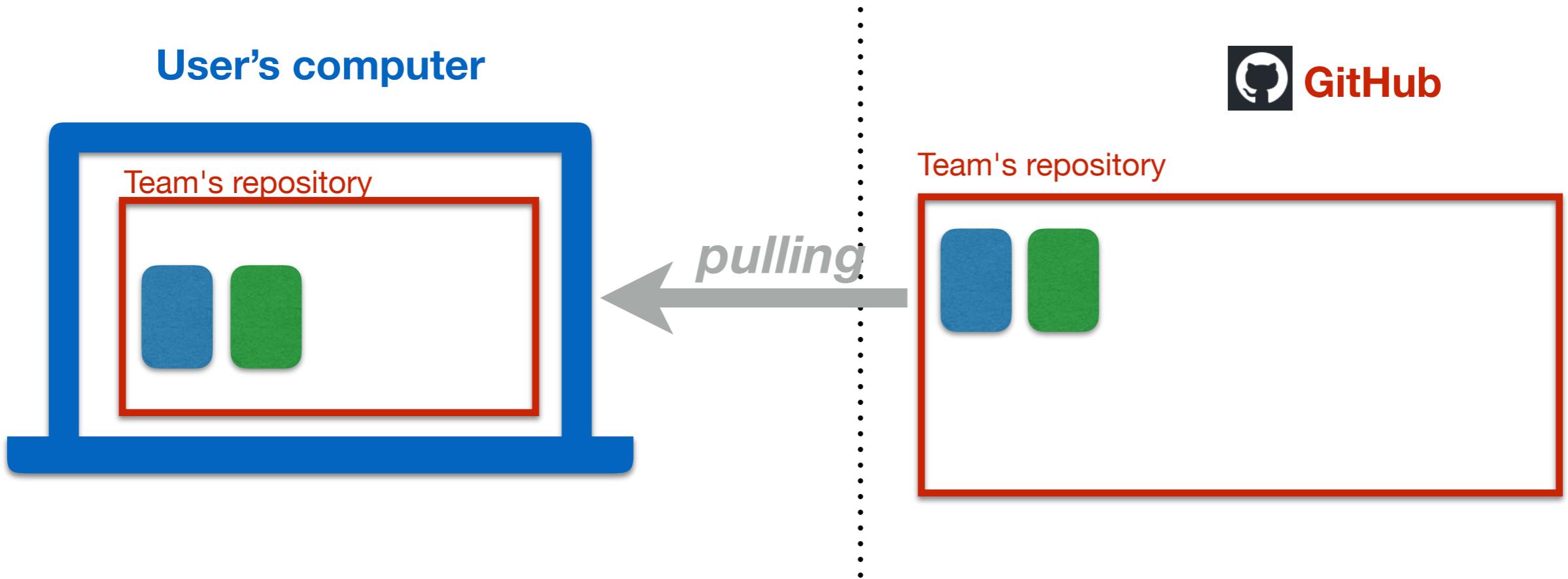
Description

Commit to master

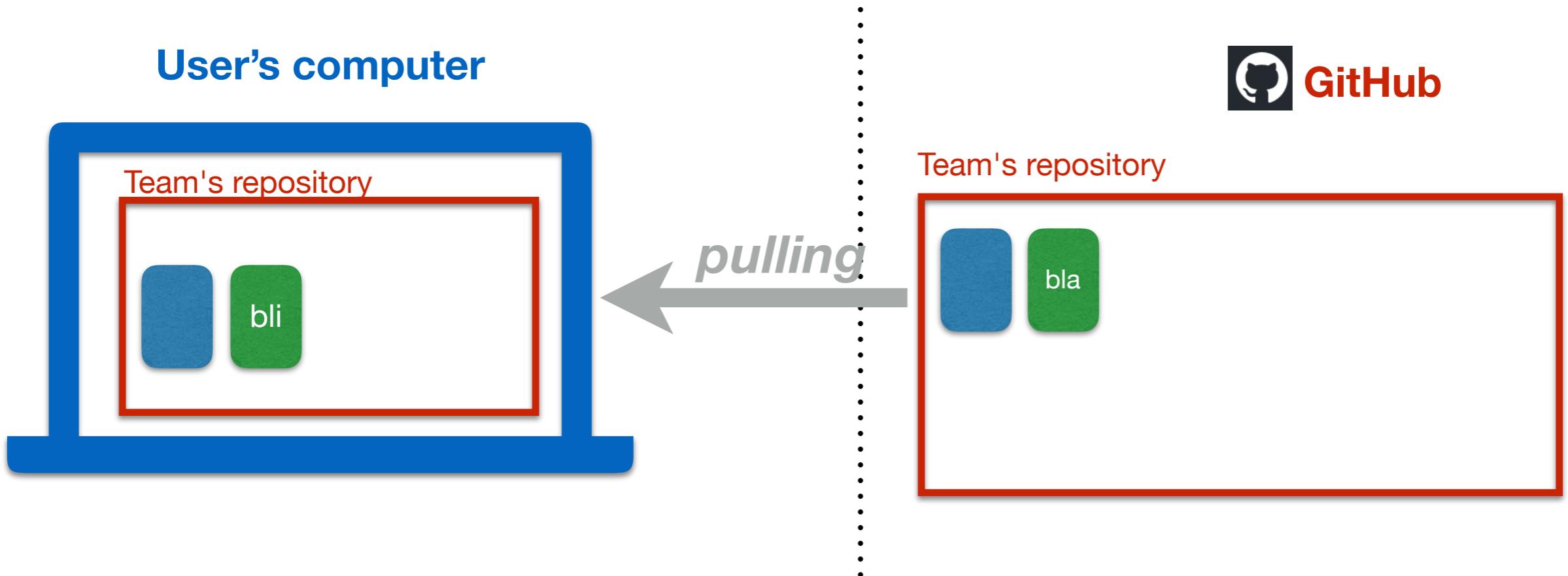
- Once the remote repository is pulled, the new file(s) are available locally



3. Pulling from the remote repository



4. conflicting changes





Current Repository **hello-world**

Current Branch **master**

Push origin
Last fetched 33 minutes ago

Changes History

0 changed files

No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

Newer Commits on Remote

! Desktop is unable to push commits to this branch because there are commits on the remote that are not present on your local branch. Fetch these new commits before pushing in order to reconcile them with your local commits.

Fetch

Summary (required)

Description

View the files in your repository in Finder
Repository menu or ⌘↑F

Show in Finder

Open the repository page on GitHub in your browser
Repository menu or ⌘↑G

View on GitHub

Commit to **master**

Committed just now
Cool modification by Jane

Undo

A modal window titled "Newer Commits on Remote" is displayed in the center. It contains a warning icon and text about remote commits. A blue "Fetch" button is at the bottom right. The background shows a GitHub desktop interface with a "Changes" tab selected, showing 0 changed files. There are also summary sections for "Description", "View in Finder", and "View on GitHub".



Current Repository **hello-world** Current Branch **master** Pull origin Last fetched just now

Changes 1 History new_file_added

1 changed file

new_file_added !

		@@ -1,3 +1,7 @@
1	1	This is new file that is added to the remote repository.
2	2	
3	3	-Awesome change by Jane! ↗↔ +<<<<< HEAD

Resolve conflicts before merging **origin/master** into **master** X

1 conflicted file

new_file_added 1 conflict Open in Visual Studio Code ▾

[Open in command line](#), your tool of choice, or close to resolve manually.

Abort merge Commit merge

hd Update new_file_added

Description

Commit to master

Committed a minute ago
Cool modification by Jane Undo

4. conflicting changes



Medizinische Fakultät Heidelberg

- Conflicting changes can be resolved with a text editor
- options depend on which editor is used

```
You, a few seconds ago | 2 authors (You and others)
This is new file that is added to the remote repository.

Accept Current Change | Accept Incoming Change | Accept Both Changes | Compare Changes
<<<<< HEAD (Current Change)
Awesome change by Jane!
=====
This is a great new modification by Joe!
>>>>> af5e9c9981b21a39ed11d09f468bea576d669191 (Incoming Change)
```

local change

changes in the
remote file

To do



Medizinische Fakultät Heidelberg

- Create your own personal GitHub account
- Register your Github user name into the Google Sheet
- all team members will be added to the corresponding GitHub repo
 - Project 03 - Team 02 → **project-03-group-02**



Medizinische Fakultät Heidelberg

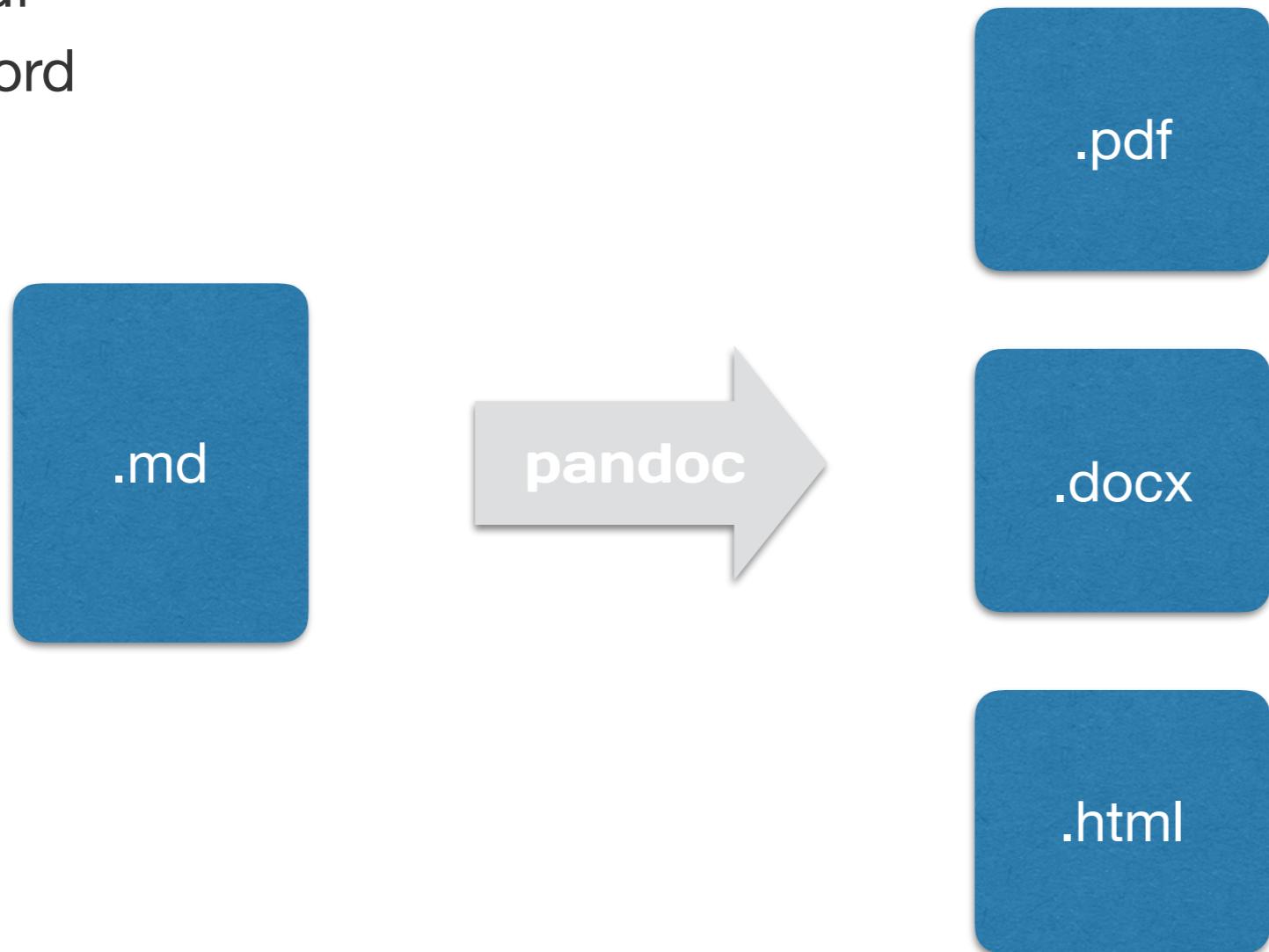
(R)markdown

Markdown



Medizinische Fakultät Heidelberg

- Markdown is a way to format plain text with a simple text editor
- Markdown documents can be converted with a **renderer** into
 - html
 - pdf
 - word



Rendering markdown



Medizinische Fakultät Heidelberg

markdown

```
# My document

## this is a header

In the text we can *highlight* or put in **bold**.

## making lists

We can make **numbered lists**:

1. first item
2. second item

or unordered lists

* first item
* second item
  + subitem
  + subitem
* third item

This is `code` which can be put inline

```bash
this is bash code
```

```python
this is python code
```

```

pdf

My document

this is a header

In the text we can *highlight* or put in **bold**.

making lists

We can make **numbered lists**:

1. first item
2. second item

or **unordered lists**

- first item
- second item
- subitem
- subitem
- third item

This is `code` which can be put inline

`this is bash code`

`this is python code`

html

My document

this is a header

In the text we can *highlight* or put in **bold**.

making lists

We can make **numbered lists**:

1. first item
2. second item

or **unordered lists**

- first item
- second item
- subitem
- subitem
- third item

This is `code` which can be put inline

`this is bash code`

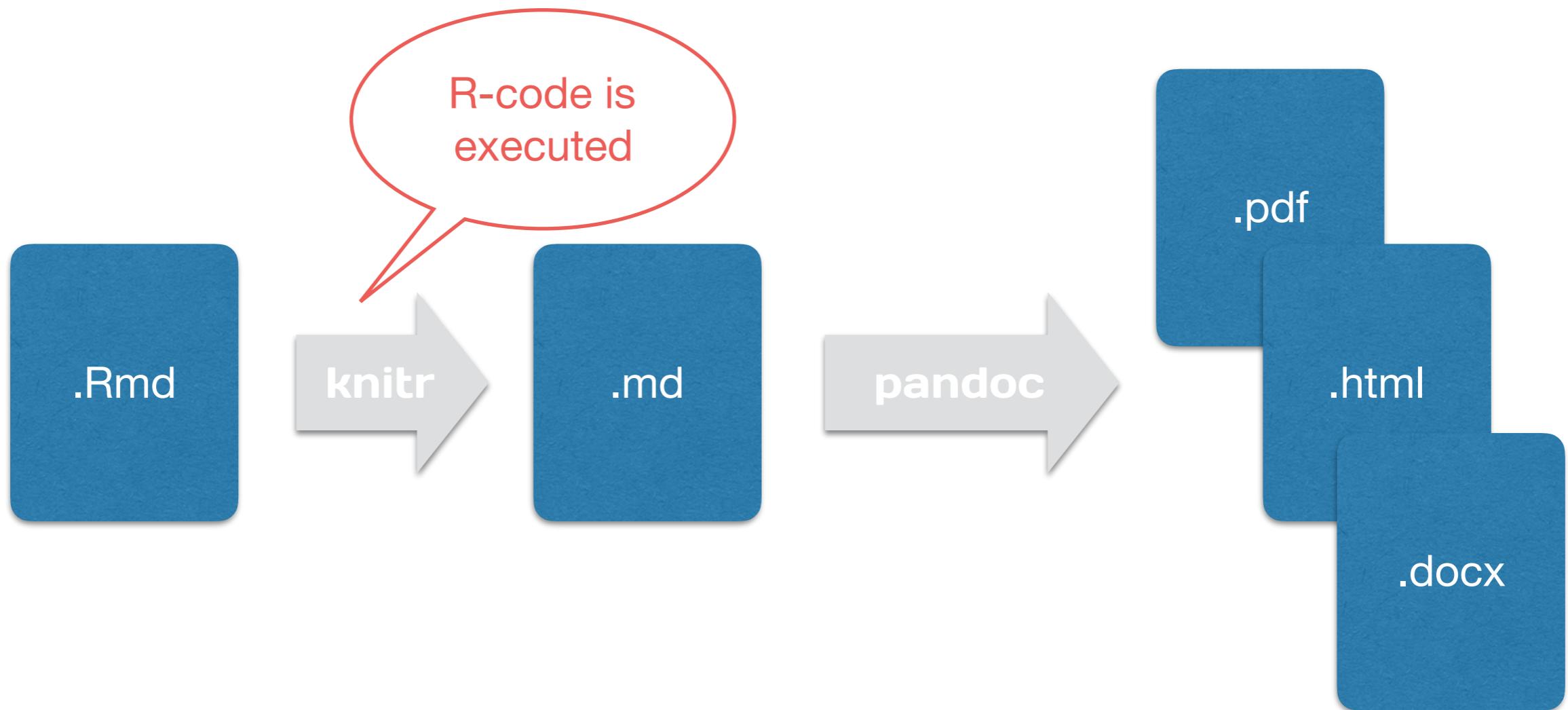
`this is python code`

Rmarkdown



Medizinische Fakultät Heidelberg

- With Rmarkdown, R-code parts can be included into the markdown document
- the R-code will be executed, the result integrated into markdown



Rmarkdown format



Medizinische Fakultät Heidelberg

```
---
```

```
title: "Project 01"
author: "Carl Herrmann"
date: "4/17/2019"
output:
  html_document:
    keep_md: yes
  pdf_document: default
---

# A Rmarkdown tutorial

This is a brief tutorial on how to use Rmarkdown to create dynamic documents

```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
knitr::opts_knit$set(root.dir='/Users/carlherrmann/Teaching/SS2019/DataAnalysis_4FS')
```

## Load the dataset

```{r read_data}
allDepMapData = readRDS('Data/depmap/DepMap19Q1_allData.RDS')
```

Now plot the distribution of the cell lines according to the tissue type

```{r plot_data}
freq = sort(table(allDepMapData$annotation$Primary.Disease))
par(las=2,mar=c(3,8,3,3));barplot(freq,horiz=TRUE, col='lightgrey')
```

```

header: set options

R code chunks

text in markdown

Markdown chunk options



Medizinische Fakultät Heidelberg

- Display options can be set for each chunk individually, or for all chunks at the beginning of the document

```
```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
knitr::opts_chunk$set(cache = TRUE)
```

valid for all chunks

- echo=TRUE : R-code is displayed in final document
- cache = TRUE : results of all chunks are cached

```
```{r plot_data,fig.height=12,fig.width=12}
freq = sort(table(allDepMapData$annotation$Primary.Disease))
par(las=2,mar=c(3,8,3,3));barplot(freq,horiz=FALSE, col='lightgrey')
````
```

valid for **this** chunks

- set height and width of output figure

# Reference



Medizinische Fakultät Heidelberg

- <https://rmarkdown.rstudio.com/>
- <https://www.rstudio.com/wp-content/uploads/2016/03/rmarkdown-cheatsheet-2.0.pdf>
-



Medizinische Fakultät Heidelberg

# Conda

# (mini)conda



Medizinische Fakultät Heidelberg

- Conda is a tool to easily install software (for example python libraries)
- Conda can be used to create environments, which contain specific software
- These environment are independent of each other, i.e. do not interfere with each other
- Example
  - one environment with Python version 2.7 + libraries
  - one environment with Python 3 + libraries